# Towards Performance Portability through an Integrated Programming Eco-System for Tensor Algebra

**Gokcen Kestor, Roberto Gioiosa, Mark Raugas**

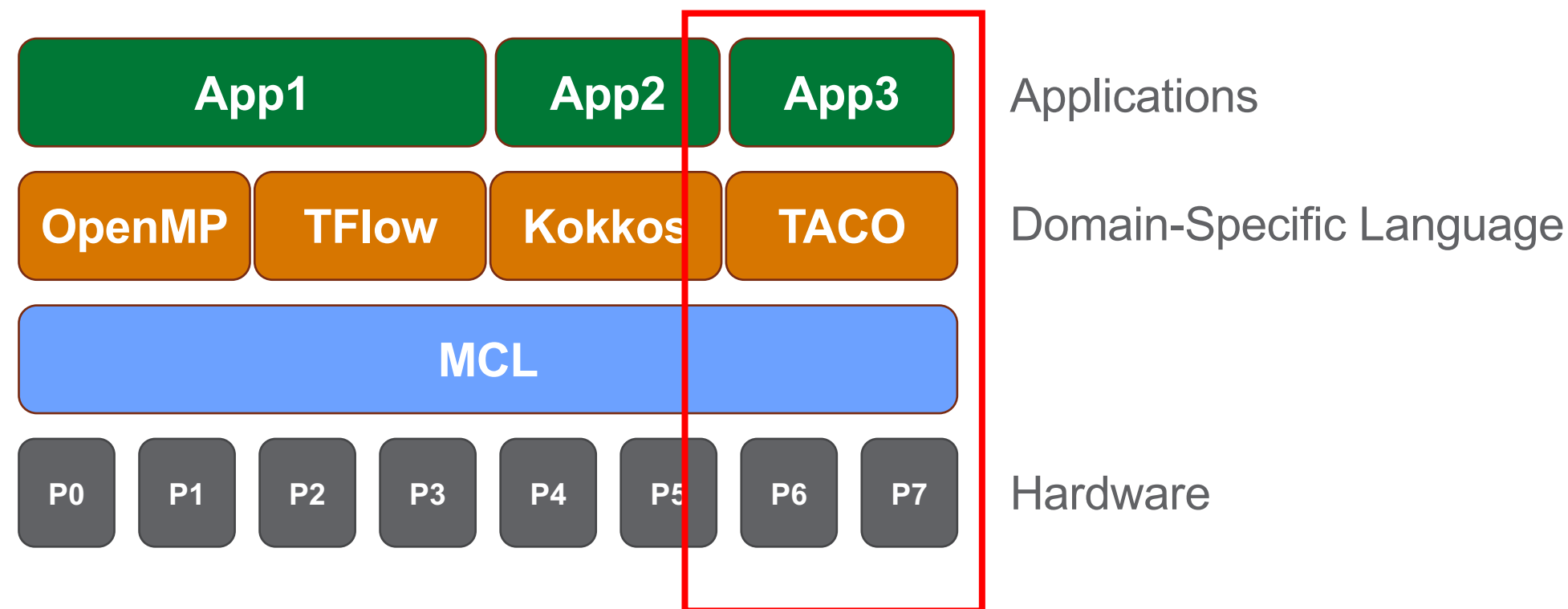Pacific Northwest National Laboratory

# Outline

- Introduction and Background
    - Tensor Algebra COmpiler (TACO)
    - Minos Computing Library (MCL)

- TACO-MCL Integrated Software Stack

- Initial Results

- Conclusions

# Motivation

- The Cambrian era is upon us:
  - Hardware landscape:
    - ✓ Many custom accelerators are being developed
    - ✓ Each HW design has its own interface, performance and energy profile
  - Software landscape:
    - ✓ Complex workflows (simulations + in-situ data analytics, simulations + AI)
    - ✓ Many programming languages and frameworks (from C/C++ to Python, TensorFlow, etc.)
- Program and performance portability has become a major concern:
  - Current HPC systems: ORNL Summit, LLNL Sierra, SNL Trinity
  - Next HPC systems: ORNL Frontier, LLNL El Capitan, ANL Aurora
- Expecting multi-device systems with several classes of devices within a single SoC (e.g., CPUs, GPUs, AI engines, FPGAs, …)
- Programming such systems is challenging!

# Proposal: A Portable Hardware/Software Stack

| App1 | App2 | App3 | Applications |
|------|------|------|--------------|
| OpenMP | TFlow | Kokkos | TACO | Domain-Specific Language |
| MCL | | | |
| P0 | P1 | P2 | P3 | P4 | P5 | P6 | P7 | Hardware |

- Scientists express their algorithm with high-level DSLs that provide domain-specific programming abstractions

- Compiler lowers DSL code to device-specific, highly-optimized code

- Dynamic runtime coordinates access to computing resources and data transfers

# Tensor Algebra COmpiler (TACO)

- TACO is a fast and versatile library for linear and tensor algebra
- C++ and Python extension to support complex tensor expression
  - Mostly focused on sparse tensor algebra*
- Automatically generate
  - Sequential CPU code
  - Parallel OpenMP code
  - NVIDIA CUDA GPU code

* Not all sparse tensor algebra operations are supported

*Fredrik Kjolstad, Shoaib Kamil, Stephen Chou, David Lugato, and Saman Amarasinghe. 2017. The tensor algebra compiler. Proc. ACM Program. Lang. 1, OOPSLA, Article 77 (October 2017), 29 pages. DOI:https://doi.org/10.1145/3133901*

```cpp
#include <iostream>
#include "taco.h"

using namespace taco;

int main(int argc, char* argv[]) {
  Format csr({Dense,Sparse});
  Format csf({Sparse,Sparse,Sparse});
  Format  sv({Sparse});

  Tensor<double> A("A", {2,3},   csr);
  Tensor<double> B("B", {2,3,4}, csf);
  Tensor<double> c("c", {4},      sv);

  // Insert data into B and c
  B(0,0,0) = 1.0;
  B(1,2,0) = 2.0;
  B(1,2,1) = 3.0;
  c(0) = 4.0;
  c(1) = 5.0;

  IndexVar i, j, k;
  A(i,j) = B(i,j,k) * c(k);

  std::cout << A << std::endl;
}
```

# TACO Example

```
1   int compute(taco_tensor_t *C, taco_tensor_t *A, taco_tensor_t *B) {
2     int C1_dimension = (int)(C->dimensions[0]);
3     int C2_dimension = (int)(C->dimensions[1]);
4     double* restrict C_vals = (double*)(C->vals);
5     int A1_dimension = (int)(A->dimensions[0]);
6     int A2_dimension = (int)(A->dimensions[1]);
7     double* restrict A_vals = (double*)(A->vals);
8     int B1_dimension = (int)(B->dimensions[0]);
9     int B2_dimension = (int)(B->dimensions[1]);
10    double* restrict B_vals = (double*)(B->vals);
11
12    #pragma omp parallel for schedule(static)
13    for (int32_t pC = 0; pC < (C1_dimension * C2_dimension); pC++) {
14      C_vals[pC] = 0.0;
15    }
16
17    #pragma omp parallel for schedule(runtime)
18    for (int32_t i0 = 0; i0 < ((A1_dimension + 31) / 32); i0++) {
19      for (int32_t i1 = 0; i1 < 32; i1++) {
20        int32_t i = i0 * 32 + i1;
21        if (i >= A1_dimension)
22          continue;
23
24        for (int32_t j = 0; j < B1_dimension; j++) {
25          int32_t jA = i * A2_dimension + j;
26          for (int32_t k = 0; k < B2_dimension; k++) {
27            int32_t kC = i * C2_dimension + k;
28            int32_t kB = j * B2_dimension + k;
29            C_vals[kC] = C_vals[kC] + A_vals[jA] * B_vals[kB];
30          }
31        }
32      }
33    }
34    return 0;
35  }
```

$$y(i) = A(i,j) * x(j)$$

CUDA code
generation for
sparse matrix-dense
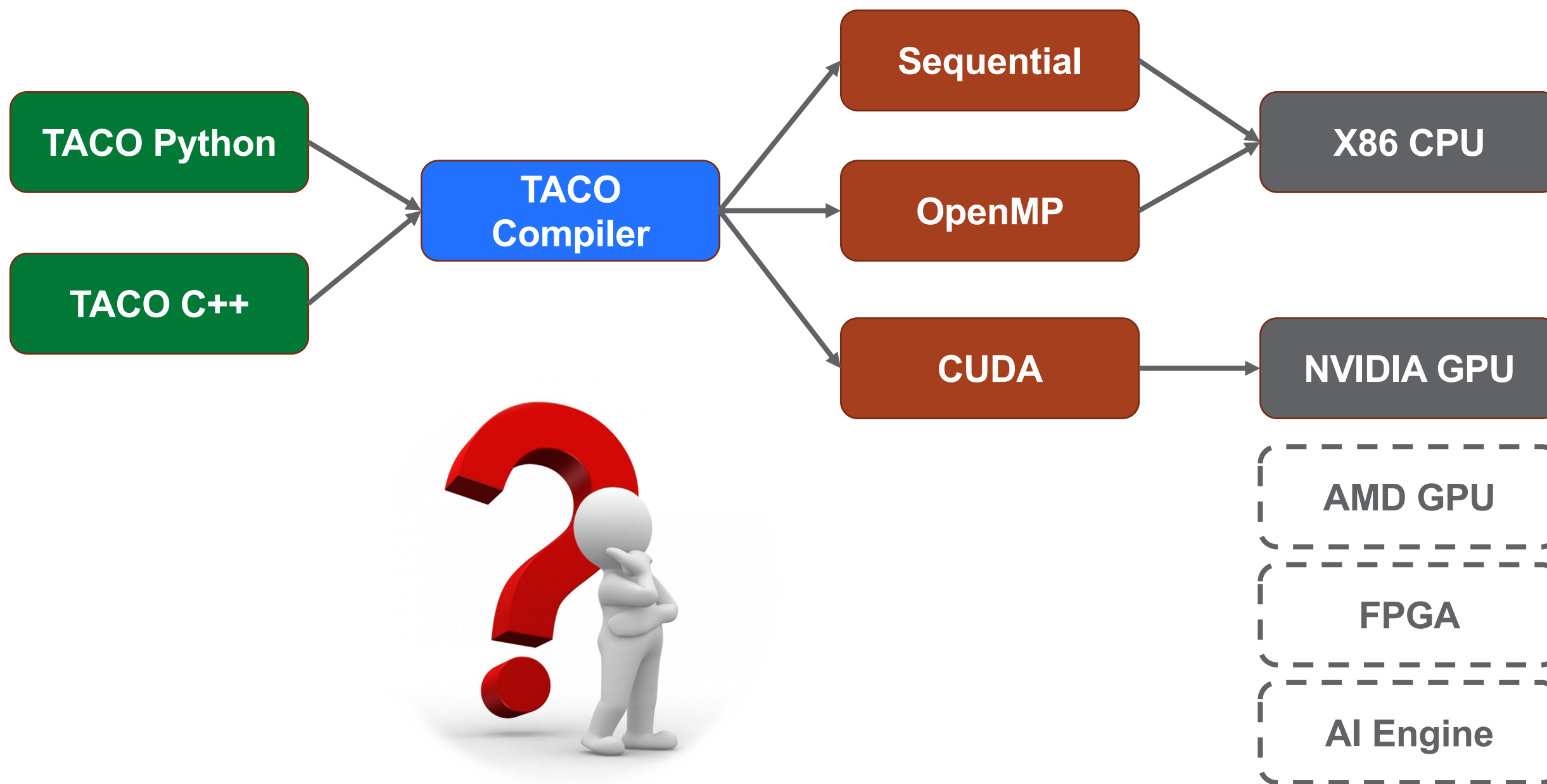vector computation

$$C(i,k) = A(i,j) * B(j,k)$$

OpenMP code
generation for
dense DGEMM
computation

```
1    __global__
2    void computeDeviceKernel0(taco_tensor_t * __restrict__ A, int32_t* i_blockStarts,
3      taco_tensor_t * __restrict__ x, taco_tensor_t * __restrict__ y){
4      int A1_dimension = (int)(A->dimensions[0]);
5      int* __restrict__ A2_pos = (int*)(A->indices[1][0]);
6      int* __restrict__ A2_crd = (int*)(A->indices[1][1]);
7      double* __restrict__ A_vals = (double*)(A->vals);
8      double* __restrict__ x_vals = (double*)(x->vals);
9      double* __restrict__ y_vals = (double*)(y->vals);
10
11     int32_t block = blockIdx.x;
12     int32_t thread = (threadIdx.x % (32));
13     int32_t warp = (threadIdx.x / 32);
14     if (threadIdx.x >= 512) {
15       return;
16     }
17
18     double workspace[7];
19     for (int32_t pworkspace = 0; pworkspace < 7; pworkspace++) {
20       workspace[pworkspace] = 0.0;
21     }
22     int32_t thr_nz = 0;
23     int32_t fpos2 = thread * 7 + thr_nz;
24     int32_t fpos1 = warp * 224 + fpos2;
25     int32_t fposA = block * 3584 + fpos1;
26     int32_t f = A2_crd[fposA];
27     if (block * 3584 + fpos1 + 7 >= A2_pos[A1_dimension]) {
28       for (int32_t thr_nz_pre = 0; thr_nz_pre < 7; thr_nz_pre++) {
29         int32_t thr_nz = thr_nz_pre;
30         int32_t fpos2 = thread * 7 + thr_nz;
31         int32_t fpos1 = warp * 224 + fpos2;
32         int32_t fposA = block * 3584 + fpos1;
33         if (fposA >= A2_pos[A1_dimension])
34           break;
35
36         int32_t f = A2_crd[fposA];
37         workspace[thr_nz_pre] = A_vals[fposA] * x_vals[f];
38       }
39     }
40     else {
41       #pragma unroll 7
42       for (int32_t thr_nz_pre = 0; thr_nz_pre < 7; thr_nz_pre++) {
43         int32_t thr_nz = thr_nz_pre;
44         int32_t fpos2 = thread * 7 + thr_nz;
45         int32_t fpos1 = warp * 224 + fpos2;
46         int32_t fposA = block * 3584 + fpos1;
47         int32_t f = A2_crd[fposA];
48         workspace[thr_nz_pre] = A_vals[fposA] * x_vals[f];
49       }
50     }
51     int32_t pA2_begin = i_blockStarts[block];
52     int32_t pA2_end = i_blockStarts[(block + 1)];
53     int32_t i_pos = taco_binarySearchBefore(A2_pos, pA2_begin, pA2_end, fposA);
54     int32_t i = i_pos;
55     for (int32_t thr_nz = 0; thr_nz < 7; thr_nz++) {
56       int32_t fpos2 = thread * 7 + thr_nz;
57       int32_t fpos1 = warp * 224 + fpos2;
58       int32_t fposA = block * 3584 + fpos1;
59       if (fposA >= A2_pos[A1_dimension])
60         break;
61
62       int32_t f = A2_crd[fposA];
63       while (fposA == A2_pos[(i_pos + 1)]) {
64         i_pos = i_pos + 1;
65         i = i_pos;
66       }
67       atomicAdd(&y_vals[i], workspace[thr_nz]);
68     }
69
70   }
```

# TACO Software Stack

# TACO-MCL: Integrated Programming Eco-System for Tensor Algebra

**TACO C++/Python Language**

**TACO-MCL Compiler**

- Automatically generate portable MCL host code and OpenCL kernels
- Break long expressions into smaller kernels for multi-device execution
- Analyze data and control flow dependencies to maximize asynchronous execution

**MCL Runtime**

**Heterogeneous Devices**

- Asynchronous task execution and overlapping of data transfers and computation
- Load balancing and resource management
- Multi-applications support

# The Minos Computing Library (MCL)

- Framework for programming extremely heterogeneous systems
  - Programming model and programming model runtime
  - **Abstract low-level architecture** details from programmers
  - **Dynamic** scheduling of work onto available resources

- Key programming features:
  - Applications factored into tasks
  - **Asynchronous** execution
  - Devices are managed by the scheduler
  - Co-schedule **independent applications**
  - Simplified APIs and programming model (based on OpenCL)

- Flexibility:
  - Scheduling framework
  - **Multiple scheduling algorithms** co-exist
  - Code portability
  - Resources allocated **at the last moment**

*Roberto Gioiosa, Burcu O. Mutlu, Seyong Lee, Jeffrey S. Vetter, Giulio Picierro, and Marco Cesati. 2020. The Minos Computing Library: efficient parallel programming for extremely heterogeneous systems. In Proceedings of GPGPU '20). ACM, New York, NY, USA, 1-10. DOI:https://doi.org/10.1145/3366428.3380770*

# Scaling Up and Down


IBM Summit


NVIDIA DGX-1

(P100/V100)


Apple iMac Pro


Apple MacBook Pro


Xilinx MPSoC ZynQ ZCU 102/106

*Same code runs on all these systems without modification*

# TACO-MCL Software Stack



By using MCL has TACO backend, TACO applications will:

- Leverage a broader classes of computing devices
- Execute in multi-device environments
- Execute in a multi-application environment
- Exploit sophisticated scheduling and load balancing algorithms

Diagram boxes:

TACO Python, TACO C++ → TACO-MCL Compiler → Sequential, OpenMP, CUDA, MCL

Sequential → X86 CPU

OpenMP → X86 CPU

CUDA → NVIDIA GPU

MCL → NVIDIA GPU, AMD GPU, FPGA, AI Engine

# TACO-MCL Workflow

```
1    #include <iostream>
2    #include "taco.h"
3
4    using namespace taco;
5
6    int main(int argc, char* argv[]) {
7      Format csr({Dense,Sparse});
8      Format csf({Sparse,Sparse,Sparse});
9      Format  sv({Sparse});
10
11     Tensor<double> A("A", {2,3},   csr);
12     Tensor<double> B("B", {2,3,4}, csf);
13     Tensor<double> c("c", {4},     sv);
14
15     // Insert data into B and c
16     B(0,0,0) = 1.0;
17     B(1,2,0) = 2.0;
18     B(1,2,1) = 3.0;
19     c(0) = 4.0;
20     c(1) = 5.0;
21
22     IndexVar i, j, k;
23     A(i,j) = B(i,j,k) * c(k);
24
25     std::cout << A << std::endl;
26   }
27
```

Original TACO application

**TACO-MCL Compiler**

C/C++ MCL driver

**MCL / CPU**

OpenCL kernel

**Computing Device**

# Experimental Results 1/2

- CCSD(1) method from NWChem
  - Coupled cluster (CC) methods are commonly used in the post Hartree-Fock ab initio quantum chemistry and in nuclear physics computation.
  - The CC workflow is composed of iterative set of excitation (singles (S), doubles (D), triples (T), and quadruples (Q)) calculations

- Testbed:
  - NVIDIA DGX-1 V100
  - 2x Intel Xeon E5-2680, 768GB memory
  - 8x NVIDIA V100, 16GM memory, NVLink
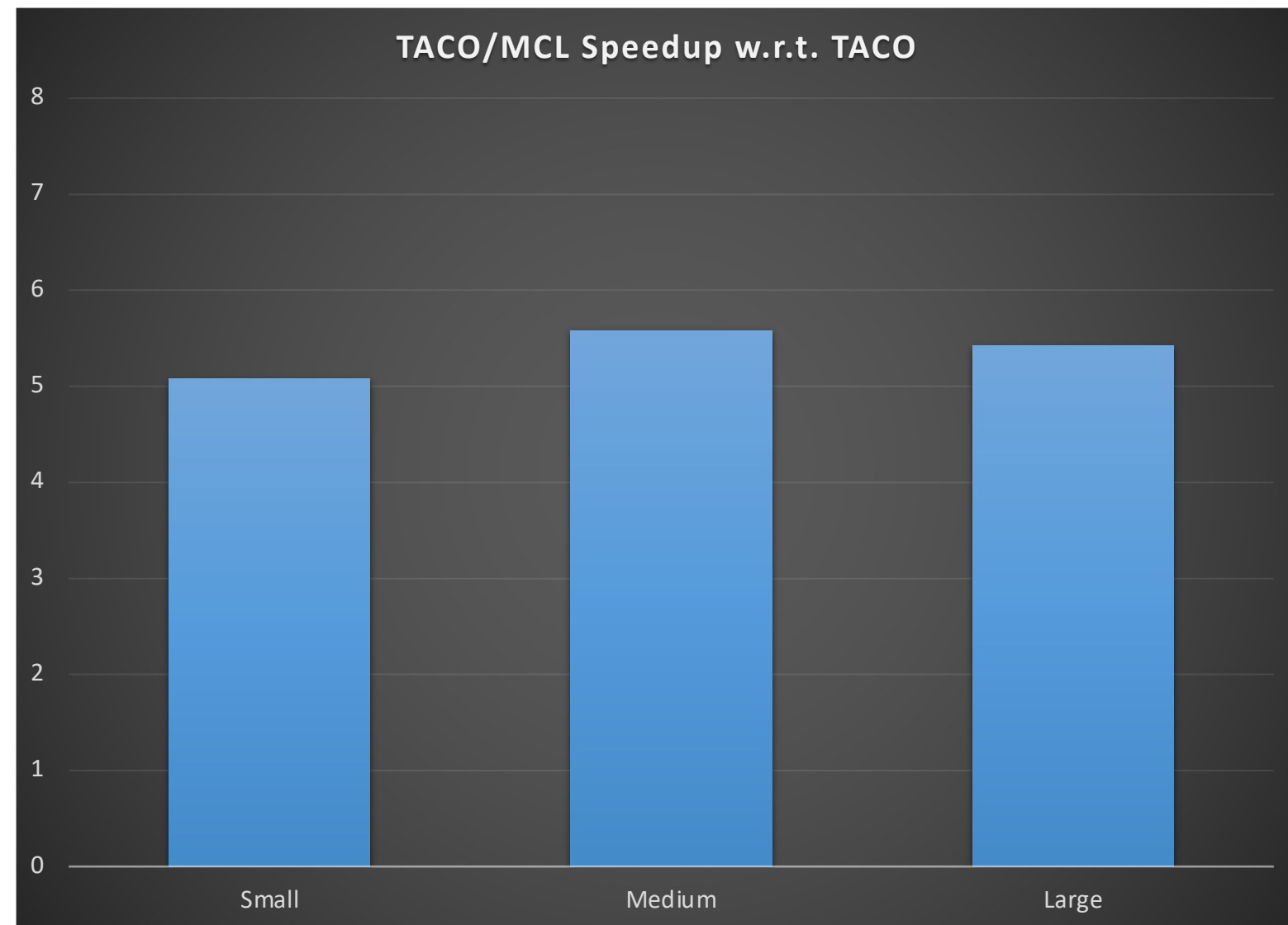
```cpp
1   #include <iostream>
2   #include "taco.h"
3   #include "utils.h"
4
5   using namespace taco;
6
7   int main(int argc, char* argv[]) {
8     if (argc != 2){
9       std::cout << "Please enter input problem size" << "\n";
10      exit(1);
11    }
12
13    int idim = atoi(argv[1]);
14
15    Format csr({Dense,Sparse});
16    Format csf({Sparse,Sparse,Sparse});
17    Format  sv({Sparse});
18
19    Format dense2d({Dense,Dense});
20    Format dense4d({Dense,Dense, Dense, Dense});
21
22    Tensor<double> i0("i0", {idim,idim}, dense2d);
23    Tensor<double> F("F", {idim, idim}, dense2d);
24    Tensor<double> V("V", {idim, idim, idim, idim}, dense4d);
25    Tensor<double> t1("t1", {idim,idim}, dense2d);
26    Tensor<double> t2("t2", {idim, idim, idim, idim}, dense4d);
27
28  // Initialization...
29
30    IndexVar i, m, n, a, e, f;
31
32    std::cout << "Computation started" << "\n";
33    i0(a, i) = F(a, i);                                              //#1
34    i0(a, i) += -2.0 * F(m, e) * t1(a, m) * t1(e, i) + F(a, e) * t1(e, i);   //#2
35    i0(a, i) += -2.0 * V(m, n, e, f) * t2(a, f, m, n) * t1(e, i);    //#3
36    i0(a, i) += -2.0 * V(m, n, e, f) * t1(a, m) * t1(f, n) * t1(e, i);  //#4
37    i0(a, i) +=  V(n, m, e, f) * t2(a, f, m, n) * t1(e, i);          //#5
38    i0(a, i) +=  V(n, m, e, f) * t1(a, m) * t1(f, n) * t1(e, i);     //%6
39    i0(a, i) += -1.0 * F(m, i) * t1(a, m);                          //#7
40    i0(a, i) += -2.0 * V(m, n, e, f) * t2(e, f, i, n) * t1(a, m);    //#8
41    i0(a, i) += -2.0 * V(m, n, e, f) * t1(e, i) * t1(f, n) * t1(a, m);  //#9
42    i0(a, i) +=  V(m, n, f, e) * t2(e, f, i, n) * t1(a, m);          //#10
43    i0(a, i) +=  V(m, n, f, e) * t1(e, i) * t1(f, n) * t1(a, m);     //#11
44    i0(a, i) +=  2.0 * F(m, e) * t2(e, a, m, i);                     //#12
45    i0(a, i) += -1.0 * F(m, e) * t2(e, a, i, m);                     //#13
46    i0(a, i) +=  F(m, e) * t1(e, i) * t1(a, m);                      //#14
47    i0(a, i) +=  4.0 * V(m, n, e, f) * t1(f, n) * t2(e, a, m, i);    //#15
48    i0(a, i) += -2.0 * V(m, n, e, f) * t1(f, n) * t2(e, a, i, m);    //#16
49    i0(a, i) +=  2.0 * V(m, n, e, f) * t1(f, n) * t1(e, i) * t1(a, m);  //#17
50    i0(a, i) += -2.0 * V(m, n, f, e) * t1(f, n) * t2(e, a, m, i);    //#18
51    i0(a, i) +=  V(m, n, f, e) * t1(f, n) * t2(e, a, i, m);          //#19
52    i0(a, i) += -1.0 * V(m, n, f, e) * t1(f, n) * t1(e, i) * t1(a, m);  //#20
53    i0(a, i) +=  2.0 * V(m, a, e, i) * t1(e, m);                     //#21
54    i0(a, i) += -1.0 * V(m, a, i, e) * t1(e, m);                     //#22
55    i0(a, i) +=  2.0 * V(m, a, e, f) * t2(e, f, m, i);              //#23
56    i0(a, i) +=  2.0 * V(m, a, e, f) * t1(e, m) * t1(f, i);         //#24
57    i0(a, i) += -1.0 * V(m, a, f, e) * t2(e, f, m, i);             //#25
58    i0(a, i) += -1.0 * V(m, a, f, e) * t1(e, m) * t1(f, i);        //#26
59    i0(a, i) += -2.0 * V(m, n, e, i) * t2(e, a, m, n);             //#27
60    i0(a, i) += -2.0 * V(m, n, e, i) * t1(e, m) * t1(a, n);        //#28
61    i0(a, i) +=  V(n, m, e, i) * t2(e, a, m, n);                   //#29
62    i0(a, i) +=  V(n, m, e, i) * t1(e, m) * t1(a, n);              //#30
63
64    i0.compile();
65    i0.assemble();
66    i0.compute();
67  }
68
```

# Experimental Results 2/2

| Problem Size | TACO (seconds) | TACO/MCL (seconds) | Speedup w.r.t. TACO |
|---|---|---|---|
| Small | 0.85 | 0.168 | 5.086 |
| Medium | 39.07 | 7.05 | 5.58 |
| Large | 1209.93 | 223.10 | 5.43 |

- TACO applications automatically scale to use all GPUs
- All problem sizes show scalability
- Expect similar speedups with larger problems
- Not ideal speedup -- WIP



TACO/MCL Speedup w.r.t. TACO

# **Conclusions**

- Program and performance portability has become a major concern

- Current and future systems feature multi-device, multi-class accelerators
  - Programming and porting applications on such systems is extremely difficult
  - Each device class has its own programming model
  - Need to manage data locality, load balancing, correctness, and resource utilization

- We developed and approach that attempts to solve the problem with an integrated software stack:
  - Users develop applications using high-level DSLs
  - Compiler lower code to targets
  - Runtime manages data locality, load balancing, and computing resources

- With TACO-MCL, original TACO applications gains
  - Access to non-NVIDIA resources (AMD/Intel GPUs, FPGAs, AI engines)
  - Transparent and automatic access to multi-device systems
  - Transparent execution in multi-applications environments (complex workflows)

# Thank you

POC:     Gokcen Kestor

gokcen.kestor@pnnl.gov